

$$y = \beta_0 + \beta_1 x + \beta_2 z$$

where x and z are indep.

$$y = \beta_0 + \beta_1 x^2 + \beta_2 z^3$$

multiple non linear regression

* Suppose X_1, \dots, X_k are independent :

random error

$$y_i = \beta_0 + \beta_1 x_1 + \dots + \beta_{k-1} x_{k-1} + \hat{\varepsilon}_i, \quad \hat{\varepsilon}_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

residual

$$\hat{y}_i = b_0 + b_1 n_{i1} + b_2 n_{i2} + \dots + b_{k-1} n_{ik-1} + \hat{e}_i, \quad n > k$$

* Least Square Method :

$$SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \sum_{i=1}^n (y_i - b_0 - b_1 n_{i1} - \dots - b_{k-1} n_{ik-1})$$

i^{th} normal equation

$$\frac{\partial L}{\partial b_0} = -2 \sum_{i=1}^n (y_i - b_0 - b_1 n_{i1} - \dots - b_{k-1} n_{ik-1}) = 0 = \sum e_i$$

$$\frac{\partial L}{\partial b_j} = -2 \sum_{i=1}^n (y_i - b_0 - b_1 n_{i1} - \dots - b_j n_{ij} - \dots - b_{k-1} n_{ik-1}) n_{ij} = 0$$

j^{th} normal equation : $\sum e_i n_{ji} = 0$

$$\left\{ \begin{array}{l} \sum_{i=1}^n (y_i - b_0 - b_1 n_{i1} - \dots - b_{k-1} n_{ik-1}) = 0 \\ \sum_{i=1}^n (y_i n_{ij} - b_0 n_{ij} - \dots - b_j n_{ij}^2 - \dots - b_{k-1} n_{ik-1} n_{ij}) = 0 \end{array} \right.$$

$$\Rightarrow \begin{cases} b_0 = \bar{Y} - b_1 \bar{X}_1 - \dots - b_K \bar{X}_K \\ b_j = \frac{\sum y_{i:nij} - b_0 n_{i:j} - \dots - b_{j-1} n_{i:j-1} - b_{j+1} n_{i:j+1} - \dots - b_{K-1} n_{i:K-1}}{\sum n_{ij}^2} \end{cases}$$

Thus, we get a system of K linear equations

- * β_0, \dots, β_K are regression coefficients
- * When estimating the parameters we only need to assume that $E(\varepsilon_i) = 0$, $\text{Var}(\varepsilon_i) = \sigma^2$ and they are valid but for testing we need to assume normality $\varepsilon_i \sim N(0, \sigma^2)$ otherwise F test fails.
- * β_i represents the change in y_i when n_i changes and the other variables are constant so β_i 's are called partial regression coefficients.
- * Least-Squares Method using Matrix Method :

$$y_i = \beta_0 + \beta_1 n_{i:1} + \beta_2 n_{i:2} + \dots + \beta_{K-1} n_{i:K-1}, \quad n > K-1$$

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}_{n \times 1} \quad \beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_{K-1} \end{pmatrix}_{K \times 1} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}_{n \times 1}$$

$$X = \begin{pmatrix} 1 & n_{1:1} & n_{1:2} & \dots & n_{1:K-1} \\ 1 & n_{2:1} & n_{2:2} & \dots & n_{2:K-1} \end{pmatrix}$$

$$Y = X \beta + \varepsilon$$

matrix form

$$\left(\begin{array}{c} \vdots \\ 1 & n_{n1} & n_{n2} & \dots & n_{nk-1} & n_{kK} \end{array} \right) \quad \text{Response } Y$$

Design Matrix
Regression
Residuals
Vector

$$E(\varepsilon) = 0 \text{ vector} \quad \text{Var}(\varepsilon) = \sigma^2 I$$

$$E(Y) = E(X\beta + \varepsilon) = X\beta + E(\varepsilon) = X\beta$$

$$\text{Var}(Y) = \text{Var}(X\beta + \varepsilon) = \text{Var}(\varepsilon) = \sigma^2$$

Estimate β by LSE :

$$\hat{Y} = X\hat{\beta} \quad Y = X\hat{\beta} + \varepsilon$$

$$\begin{aligned}
 L &= \sum_{i=1}^n e_i^2 = (e_1, \dots, e_n) \begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix} = e^T \cdot e \\
 &= (Y - \hat{Y})^T (Y - \hat{Y}) \\
 &= (Y - X\hat{\beta})^T (Y - X\hat{\beta}) \\
 &= (Y^T - \hat{\beta}^T X^T) (Y - X\hat{\beta}) \\
 &= Y^T Y - Y^T X \hat{\beta} - \hat{\beta}^T X^T Y - \hat{\beta}^T X^T X \hat{\beta} \\
 \text{as } c^T &= c \text{ where } = Y^T Y - (\hat{\beta}^T X^T Y)^T - \hat{\beta}^T X^T Y - \hat{\beta}^T X^T X \hat{\beta} \\
 c \text{ is a scalar} &\Leftarrow = Y^T Y - 2\hat{\beta}^T X^T Y^T - \hat{\beta}^T X^T X \hat{\beta}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial}{\partial \hat{\beta}} &= -2X^T Y - 2X^T X \hat{\beta} = 0 \Rightarrow X^T X \hat{\beta} = X^T Y \\
 &\Rightarrow \hat{\beta} = (X^T X)^{-1} X^T Y
 \end{aligned}$$

$$\text{Hence, } \hat{\beta}_{OLS} = (X^T X)^{-1} X^T Y$$

When $K=1$, we get back to simple linear regression.

- * In the chocolate company example, some other regressors

could be the number of sales person, etc. In multiple linear regression we have $k-1$ regressors; x_{11}, \dots, x_{k-1}

i	x_{i1}	$x_{i2} \dots$	x_{ik-1}
1	x_{11}	$x_{12} \dots$	x_{1k-1}
2	x_{21}	$x_{22} \dots$	x_{2k-1}
\vdots	\vdots	\vdots	\vdots
n	x_{n1}	$x_{n2} \dots$	x_{nk-1}

$\det(x^T x) = 0$ not invertible
and the columns are linearly dependent otherwise they are independent and $\det(x^T x) \neq 0$
so some features might be redundant

$$\begin{aligned}\mathbb{E}(\hat{\beta}) &= \mathbb{E}((x^T x)^{-1} x^T y) \\ &= (x^T x)^{-1} x^T \mathbb{E}(x\beta + \varepsilon) \\ &= (x^T x)^{-1} x^T (x\beta + \mathbb{E}(\varepsilon)) , \mathbb{E}(\varepsilon) = 0 \\ &= (x^T x)^{-1} (x^T x) \beta \\ &= \beta \text{ unbiased estimator}\end{aligned}$$

$$x^T = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_{11} & x_{21} & \dots & x_{n1} \\ x_{12} & x_{22} & \dots & x_{n2} \\ \vdots & \vdots & & \vdots \\ x_{1k-1} & x_{2k-1} & \dots & x_{nk-1} \end{pmatrix}_{k \times n}$$

$x^T x$ is symmetric, we call it the variance-covariance matrix.

* Gauss Markov Theorem :

The OLS (ordinary least square) estimates $\hat{\beta}$ at β is the best linear unbiased estimator (BLUE).

* Proof : Elements of Statistical Learning

Linear Regression Analysis by Montgomery

* Def: BLUE

An estimator $\hat{\theta}$ is called BLUE for θ if

- 1) $\hat{\theta}$ is a linear combination of samples obs ($\hat{\beta}_i = \sum c_i y_i$ in linear)
- 2) $\text{Var}(\hat{\theta}) \leq \text{Var}(\theta')$ where $\hat{\theta}$ and θ' are two unbiased estimators.

$$\begin{aligned}\text{Var}(\hat{\beta}) &= \text{Var}(c(x^T x)^{-1} x^T y) \\ &= ((x^T x)^{-1} x^T) ((x^T x)^{-1} x^T)^T \text{Var}(y) \\ &= (x^T x)^{-1} x^T x (x^T x)^{-1} \sigma^2 I \\ &= I \cdot (x^T x)^{-1} \sigma^2 I \\ &= \sigma^2 (x^T x)^{-1} \\ &= \text{MSRes} (x^T x)^{-1} \quad \text{when } \sigma \text{ unknown, } n-k \text{ deg.}\end{aligned}$$

$$\text{SSRes} = \sum e_i^2 \sim N(0, \sigma^2) \Rightarrow \frac{e_i^2}{\sigma^2} \sim N(0, 1)$$

$$\Rightarrow \frac{e_i^2}{\sigma^2} \sim \chi^2_1$$

because we have K normal eq.

$$\Rightarrow \frac{\sum e_i^2}{\sigma^2} \sim \chi^2_{n-k}$$

$$E(\text{MSRes}) = \sigma^2$$

$$SS_T = \sum_{i=1}^n (y_i - \bar{y})^2 \quad n-1 \text{ df as we have one restriction } \sum y_i - \bar{y} = 0$$

$$\text{MSRes} = \frac{\text{SSRes}}{n-k}$$

Then, SSRes has $k-1$ degrees of freedom

* Example: In R.

$$1) Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2, \text{ compute } \beta_i \text{'s } i=0,1,2$$

2) FInd Anova table

Source of variability	DF	SS	MS	F-value
Regression	$k-1$	$\sum (\hat{Y} - \bar{Y})^2$	$SS_{\text{reg}} / (k-1)$	$F = \frac{MS_{\text{reg}}}{MS_{\text{res}}} \sim F_{k-1, n-k}$
Residual	$n-k$	$\sum (Y_i - \hat{Y})^2$	$SS_{\text{res}} / (n-k)$	
Total	$n-1$	$\sum (Y_i - \bar{Y})^2$	$SS_T / (n-1)$	

$$F_{\text{obs}} > F_{1, n-k | 0.05}$$

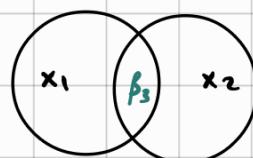
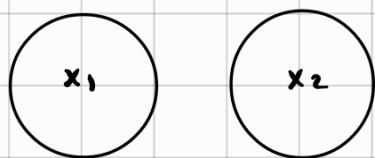
$H_0: \beta_1 = \beta_2 = 0$ rejected
otherwise MLR is not useful

3) Calculate R^2

$$R^2 = \frac{SS_{\text{reg}}}{SS_T}$$

4) Calculate the estimated variance of $\hat{\beta}$.

* Interactions :



$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \varepsilon$$

* Global Test :

$$H_0: \beta_j = 0 \quad j = 1(1)k-1 \quad \text{vs} \quad H_1: \beta_j \neq 0 \quad \exists j \in \{0, \dots, k-1\}$$

$$\frac{SS_{\text{Res}}}{\sigma^2} \sim \chi^2_{n-k}$$

$$\frac{MS_{\text{Reg}}}{MS_{\text{Res}}} \sim F_{k-1, n-k}$$

$$\frac{SS_{\text{Reg}}}{\sigma^2} \sim \chi^2_{k-1}$$

reject H_0 if $F_{\text{cal}} > F_{\alpha, k-1, n-k}$ i.e. $\exists \beta_j \neq 0$

* Partial Test : Which regressor is significant?

What is the random response to dependent variable?

$$H_0: \beta_j = 0 \quad \text{vs} \quad H_1: \beta_j \neq 0 \quad j = 1(1)K-1$$

$$\hat{\beta}_j \sim N(\beta_j, \sigma^2 (X^\top X)^{-1}_{jj}) \Rightarrow \frac{\hat{\beta}_j - \beta_j}{\sqrt{\sigma^2 (X^\top X)^{-1}_{jj}}} \sim N(0, 1)$$

under H_0 : $\frac{\hat{\beta}_j}{\sqrt{MSRes(X^\top X)^{-1}_{jj}}} \sim t_{\alpha/2, n-K}$

If $|t| > t_{\alpha/2, n-K}$, we reject H_0 in this case we do not need to do a global test as $\exists \beta_j \neq 0$

* Confidence Interval: $\frac{\hat{\beta}_j - \beta_j}{\sqrt{MSRes(X^\top X)^{-1}_{jj}}} \sim t_{\alpha/2, n-K}$

$$P\left(\left|\frac{\hat{\beta}_j - \beta_j}{\sqrt{MSRes(X^\top X)^{-1}_{jj}}}\right| < t_{\alpha/2, n-K}\right) = 1 - \alpha$$

$$P\left(\beta_j - t_{\alpha/2, n-K} \sqrt{MSRes(X^\top X)^{-1}_{jj}} < \hat{\beta}_j < \beta_j + t_{\alpha/2, n-K} \sqrt{MSRes(X^\top X)^{-1}_{jj}}\right) = 1 - \alpha$$

* Proof of Gauss Markov Theorem:

We want to show the OLS is the best linear unbiased estimator at β .

By OLS: $\hat{\beta} = (X^\top X)^{-1} X^\top Y$

$$E(\hat{\beta}) = \beta \quad \text{Var}(\hat{\beta}) = \sigma^2 (X^\top X)^{-1}$$

Now, consider another unbiased estimator of β , call it $\tilde{\beta}$. We use the fact that any other estimator of β can be written as

$$\tilde{\beta} = ((X^T X)^{-1} X^T + B) Y + \beta_0, \quad B \in \mathcal{M}_{k \times n}(\mathbb{R})$$

$$\beta_0 \in \mathcal{M}_{n \times 1}(\mathbb{R})$$

$$\begin{aligned} E(\tilde{\beta}) &= E((X^T X)^{-1} X^T Y) + E(\beta Y) + E(\beta_0) \\ &= \beta + B E(X\beta + \epsilon) + \beta_0 \\ &= \beta + B X \beta + \beta_0, \quad \beta \in \mathcal{M}_{k \times n}(\mathbb{R}) \end{aligned}$$

Note that $B X$ and β_0 are zero as we assumed that $\tilde{\beta}$ is unbiased

$$\begin{aligned} \text{Var}(\tilde{\beta}) &= \text{Var}((X^T X^{-1}) X^T + B) Y + \beta_0 \\ &= \text{Var}((X^T X^{-1}) X^T + B) Y \\ &= ((X^T X^{-1}) X^T + B)((X^T X^{-1}) X^T + B)^T \text{Var}(Y) \\ &= ((X^T X)^{-1} X^T + B)(X(X^T X)^{-1} + B^T) \sigma^2 \\ &= \sigma^2 ((X^T X)^{-1} + (X^T X)^{-1} (B X)^T + (B X) (X^T X)^{-1} + \beta \beta^T) \\ &= \sigma^2 ((X^T X)^{-1} + \beta \beta^T) \end{aligned}$$

Therefore, $\text{Var}(\tilde{\beta}) > \text{Var}(\hat{\beta})$ and thus OLS estimator of β is the BLUE because $\beta \beta^T > 0$